

# Searching for gravitational-wave transients with a qualitative signal model: seedless clustering strategies

Eric Thrane<sup>1, a</sup> and Michael Coughlin<sup>2</sup>

<sup>1</sup>*LIGO Laboratory, California Institute of Technology, MS 100-36, Pasadena, CA, 91125, USA*

<sup>2</sup>*University of Cambridge, Cambridge, CB2 1TN, United Kingdom*

(Dated: September 5, 2013)

Gravitational-wave bursts are observable as bright clusters of pixels in spectrograms of strain power. Clustering algorithms can be used to identify candidate gravitational-wave events. Clusters are often identified by grouping together seed pixels in which the power exceeds some threshold. If the gravitational-wave signal is long-lived, however, the excess power may be spread out over many pixels, none of which are bright enough to become seeds. Without seeds, the problem of detection through clustering becomes more complicated. In this paper we investigate seedless clustering algorithms in searches for long-lived narrowband gravitational-wave bursts. Using four astrophysically motivated test waveforms, we compare a seedless clustering algorithm to two algorithms using seeds. We find that the seedless algorithm can detect gravitational-wave signals (at fixed false-alarm and false-dismissal rate) at distances between 150–200% greater than those achieved with the seed-based clustering algorithms, corresponding to significantly increased detection volumes: 420–740%. This improvement in sensitivity may extend the reach of second-generation detectors such as Advanced LIGO and Advanced Virgo deeper into astrophysically interesting distances.

PACS numbers: 95.75.-z, 04.30.-w

## I. INTRODUCTION

Searches for gravitational-wave (GW) transients typically fall into two classes. “Burst” searches employ only minimal assumptions to target unmodeled or difficult-to-model GW sources. Other GW sources, such as coalescing neutron stars / black holes, produce readily predictable waveforms, making it possible to carry out a near-optimal search with a matched filter template bank. However, it is also possible to design a GW transient search in between these two opposite ends of the spectra, where some information about the signal model is known, but not enough to produce a reliable template bank. In this paper we investigate the possibility of GW transient searches for which we have a qualitative signal model, focusing in particular on models predicting GW signals, which are long-lived  $\gtrsim 10$  s and narrowband, but which are otherwise poorly constrained.

Long-lived narrowband GW transients have been proposed to originate in a variety of astrophysical processes, most notably, in newborn neutron stars [1–4] and black hole accretion disks following stellar collapse [5–7]. Long-lived GW transients can be observed with excess strain power algorithms [1, 8]. Signals show up as curved tracks on  $ft$ -maps (spectrograms) of strain power, see Fig.1.

A number of clustering algorithms have been proposed to identify statistically significant GW signatures in strain power spectrograms, see, e.g., [8–13]. Most existing algorithms rely on the use of seeds: spectrogram pixels with excess power above some threshold [25]. The idea behind seed-based algorithms is that sufficiently

loud GW signals induce excess power, which leads to the creation of seeds along a spectrogram track. The clustering algorithm connects neighboring seeds in order to form a cluster. (Different algorithms use different rules for connecting seeds.) Next, the clustered seeds are combined to produce a detection statistic, which is used to determine if the cluster is consistent with detector noise.

One of the advantages of seed-based clustering is that only minimal assumptions need be made about the signal. While different clustering rules may be better or worse for different signal models—e.g., narrowband tracks versus broadband blobs—most seed-based clustering algorithms can effectively cluster signals with arbitrary spectrographic morphology given a sufficiently high signal-to-noise ratio.

One disadvantage of seed-based clustering is that the signal must be loud enough to create seeds in the first place or the whole enterprise is doomed. As we seek to study longer and weaker GW signals, this becomes increasingly problematic. For a fixed energy budget, the average excess power in each of  $N$  spectrogram pixels scales like  $1/N$ . In other words, long signals are less likely to induce seeds than short signals, all else equal.

Here we investigate seedless clustering algorithms designed to target long and weak signals. We propose a seedless clustering algorithm that will enforce additional assumptions about the signal model: that it is long-lived and narrowband. By making these assumptions, we sacrifice some of the flexibility of seed-based clustering algorithms for improved sensitivity to a specific class of signal models.

The remainder of this paper is organized as follows. In Section II we formulate the problem of detecting long-lived narrowband GW transients as a pattern recognition problem: how to detect tracks from GWs in strain power

---

<sup>a</sup>Electronic address: ethrane@ligo.caltech.edu

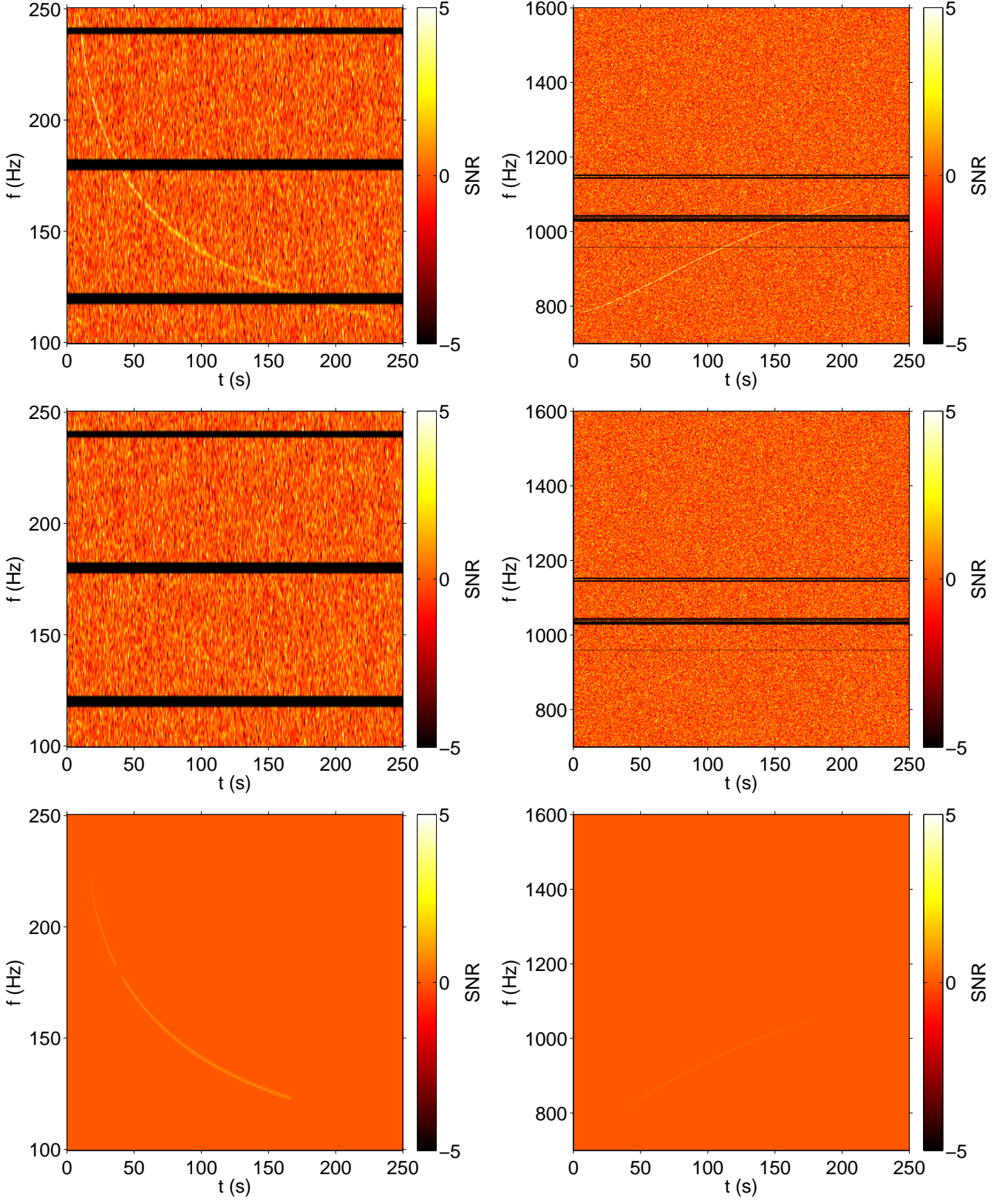


FIG. 1: Injection recovery with seedless clustering using simulated Advanced LIGO noise. Top row: SNR spectrograms for relatively nearby signals. Left is a  $d = 150$  Mpc accretion disk instability signal (ADI 2) and right is a  $d = 16$  Mpc fallback accretion signal (FA 2); see Tab. I. The black horizontal lines are notches due to instrumental artifacts. Second row: the same as the first row, but the injected signals are further away,  $d = 360$  Mpc (left) and  $39$  Mpc (right), and so the SNRs are less by  $\approx 6\times$ . The tracks are all but invisible to the naked eye. Bottom row: the loudest recovered tracks obtained by analyzing the second-row spectrograms with stochtrack ( $T = 2 \times 10^8$  trials). Both clusters have FAP  $< 0.1\%$ .

spectrograms. In Section III we describe clustering algorithms with seeds and introduce an alternative seedless clustering algorithm. In Section IV, we describe a Monte Carlo study comparing the sensitivity of seed-based and seedless clustering algorithms. We also repeat the analysis using recolored initial LIGO noise (with an unphysical time shift) to study whether the Monte Carlo results hold for non-ideal detector noise. Finally, in Section V, we summarize our findings, describe the limitations of our proposed algorithm, and discuss possibilities for future work.

## II. FORMALISM

Searches for long GW transients can be cast as pattern recognition problems [8]. Strain time series  $s_I(t')$  for detector  $I$  are divided into segments of duration  $\delta t$  with start times  $t$ . (Note that  $t'$  denotes sampling times whereas  $t$  denotes segment start times.) The Fourier transform of the detector- $I$  strain data in segment  $t$  is denoted  $\tilde{s}_I(t; f)$ .

Following [8], we define an estimator for strain cross-power in the  $IJ$  detector pair  $\hat{Y}_{IJ}(t; f)$  with associated variance  $\hat{\sigma}^2(t; f)$ :

$$\begin{aligned}\hat{Y}(t; f) &= \frac{2}{\mathcal{N}} \text{Re} \left[ Q_{IJ}(t; f | \hat{\Omega}) \tilde{s}_I^*(t; f) \tilde{s}_J(t; f) \right] \\ \hat{\sigma}^2(t; f) &= \frac{1}{2} \left| Q_{IJ}(t; f | \hat{\Omega}) \right|^2 P'_I(t; f) P'_J(t; f).\end{aligned}\quad (1)$$

Here  $\mathcal{N}$  is an FFT normalization factor and  $Q_{IJ}(t; f | \hat{\Omega})$  is a filter function, which takes into account the relative time delays and the  $IJ$  detector responses for a source located in the direction of  $\hat{\Omega}$ . The filter function is defined such that  $\hat{Y}(t; f)$  is an unbiased estimator for GW power [8]. Meanwhile,  $P'_I(t; f)$  and  $P'_J(t; f)$  are the auto-power spectral densities for detectors  $I$  and  $J$  in the segments neighboring  $t$ . For additional details, see [8].

We define signal-to-noise ratio for a spectrogram pixel at  $(t; f)$  as

$$\rho(t; f) \equiv \hat{Y}(t; f) / \hat{\sigma}(t; f). \quad (2)$$

An array of  $\rho(t; f)$  can be visualized as an  $ft$ -map as in Fig. 1. Detector noise is distributed quasi-normally with mean  $\langle \rho(t; f) \rangle = 0$  while GW signals produce positive contributions to  $\rho(t; f)$ . A loud long-lived narrowband transient, therefore, appears as a track of bright pixels in a spectrogram of  $\rho(t; f)$ . If the GW signal is very weak, the track may not be visible by eye, though, there is still a statistical excess in  $\rho(t; f)$  along the GW track.

The job of a clustering algorithm is to identify a cluster of pixels  $\Gamma$ , which, subject to some set of clustering rules, is more likely than any other cluster to be associated with a GW signal. In order to determine which cluster among many is loudest, and in order to determine the statistical significance of a cluster, it is necessary to define a detection statistic characterizing the loudness of the

entire cluster. Following [8], we define the cluster signal-to-noise ratio as

$$\text{SNR}_{\text{tot}} \equiv \frac{\sum_{\{t; f\} \in \Gamma} w(t; f) Y(t; f)}{\left( \sum_{\{t; f\} \in \Gamma} w^2(t; f) \sigma^2(t; f) \right)^{1/2}}. \quad (3)$$

Here  $w(t; f)$  is a weight factor, which can be chosen to emphasize certain frequencies and times depending on the detector noise, the expected GW signal, both, or neither.

By performing many pseudo experiments with Monte Carlo or time-shifted detector noise, it is possible to measure the probability density function  $p(\text{SNR}_{\text{tot}})$  from which we determine the threshold  $\text{SNR}_{\text{tot}}^{\text{th}}$  required for a detection at fixed false-alarm probability (FAP):

$$\int_0^{\text{SNR}_{\text{tot}}^{\text{th}}} d(\text{SNR}_{\text{tot}}) p(\text{SNR}_{\text{tot}}) = 1 - \text{FAP}. \quad (4)$$

The sensitivity of a clustering algorithm to a specific source can be characterized by the distance to which it can detect the source with  $\text{SNR}_{\text{tot}} \geq \text{SNR}_{\text{tot}}^{\text{th}}$  with fixed FAP and fixed false-dismissal probability (FDP). In this paper we define detection distance  $d_0$  as the distance at which a GW signal can be observed above threshold with FAP = 0.1% and FDP = 50%. Detection distance is always defined for a specific gravitational waveform (model), so below we present results for several models.

## III. CLUSTERING

In this section we discuss how different clustering algorithms can be used to identify tracks of excess power in spectrograms of  $\rho(t; f)$ .

### A. Seed-based clustering

The first step for any seed-based algorithm is to apply a threshold in order to identify seeds:

$$\rho(t; f) > \rho_{\text{th}}. \quad (5)$$

The threshold is a tunable parameter that can be chosen so as to maximize  $d_0$ . If  $\rho_{\text{th}}$  is too small, there will be many seeds due to noise fluctuations, which leads to many loud noise clusters, ultimately harming the sensitivity of the search. In fact, if  $\rho_{\text{th}}$  is made sufficiently small, the typical density of seed pixels will be so great that seeds from noise fluctuations will form a single large cluster spreading throughout the spectrogram. On the other hand, if  $\rho_{\text{th}}$  is too large, only very loud signals will create seeds. We find empirically that  $\rho_{\text{th}} \approx 0.75$  maximizes  $d_0$  for the seed-based clustering algorithms considered here.

Next, the seeds are combined to produce clusters. There are myriad ways of clustering seeds. Linear clustering algorithms (e.g., [11]) combine seeds that fall within



a fixed distance of each other. Density-based clustering algorithms (e.g., [9]) require that the number of seeds per unit area exceeds some threshold in order to be joined. The “locust” algorithm [10], meanwhile, is a local wandering scheme in which the two most significant neighboring seeds in some box are connected iteratively until no more seeds are available to connect. It is also possible to combine the seeds along predefined paths specified by polynomials using a Hough algorithm [10]. In the comparison that follows, we employ a linear clustering algorithm [11] and a density based algorithm [9], both of which are in use in GW transient analyses [14, 15].

One advantage of seed-based clustering is that most implementations, as a rule of thumb, can be made to operate with relatively modest computational resources. Reducing a large number of pixels in a  $\rho(t; f)$  spectrogram to a handful of seeds simplifies the clustering problem.

One disadvantage of seed-based clustering is that the excess strain power from long signals is spread out over many pixels and may therefore fail to produce seeds. Another disadvantage arises from the presence of instrumental noise lines present GW strain data; see Fig. 1. Noise lines must be notched to avoid numerous clusters from non-stationary noise. The notches, in turn, create gaps over which it may be difficult to join seeds. In the next subsection, we show how seedless clustering can overcome both of these obstacles.

## B. Seedless clustering

A seedless clustering algorithm does not apply a threshold to  $\rho(t; f)$ . An example of a previously proposed seedless clustering algorithm is the Radon algorithm [8], which integrates  $\rho(t; f)$  along every possible straight line that can be drawn through  $\rho(t; f)$ . There are a number of limitations associated with the Radon algorithm, which we pause to study in order so that we might illuminate the path to a more effective clustering strategy.

First, the Radon algorithm assumes the track is well-described as a straight line in  $ft$ -space, which is a poor approximation for many realistic signals, see Fig. 1. Second, it assumes that the signal persists for the duration of the spectrogram (or until the line intersects the top/bottom edges). Finally, the background is needlessly increased by including nearly vertical lines, corresponding to short times, which do not conform to the assumed long-lived signal model.

We endeavor to address these shortcomings with a new seedless algorithm, which we call *stochtrack*. The basic idea of *stochtrack* is to integrate  $\rho(t; f)$  along monotonic  $f(t)$  curves with arbitrary start and stop times subject to the constraint that the total duration is at least  $t_{\min}$  taken here to be 20–100 s depending on the model. By allowing for curved tracks, we aim to better fit plausible GW signals.

The algorithm works as follows:

1. Choose a random triplet of start-time, mid-time, and stop-time ( $t_{\text{start}}, t_{\text{mid}}, t_{\text{stop}}$ ) such that  $(t_{\text{stop}} - t_{\text{start}}) \geq t_{\min}$  and  $t_{\text{start}} < t_{\text{mid}} < t_{\text{end}}$ .
2. Choose a random triplet of start-frequency, mid-frequency, and stop-frequency ( $f_{\text{start}}, f_{\text{mid}}, f_{\text{stop}}$ ) such that  $f_{\text{start}} \leq f_{\text{mid}} \leq f_{\text{end}}$  (up-chirping) or  $f_{\text{start}} \geq f_{\text{mid}} \geq f_{\text{end}}$  (down-chirping).
3. These two triplets correspond to three ordered pairs of  $(f, t)$ . Using the three ordered pairs as control points, form a quadratic Bézier curve [16] denoted  $\Gamma$ . (Other curve parameterizations, such as a cubic spline, are possible as well.)
4. Following Eq. 3, perform a weighted sum of the values of  $\rho(t; f)$  in  $\Gamma$  to calculate  $\text{SNR}_{\text{tot}}$ .
5. Repeat the previous steps  $T$  times. Record the cluster with the largest value of  $\text{SNR}_{\text{tot}}$ .

Above we have described the *stochtrack* algorithm in terms of a for-loop, but in practice it can be more computationally efficient to work with  $T$ -dimensional vectors of ordered pairs:  $(\vec{t}_{\text{start}}, \vec{f}_{\text{start}})$ ,  $(\vec{t}_{\text{mid}}, \vec{f}_{\text{mid}})$ , and  $(\vec{t}_{\text{end}}, \vec{f}_{\text{end}})$ .

In order to explore some of the computational subtleties of this calculation, it is worthwhile to consider a concrete example. Consider a 151 Hz  $\times$  250 s spectrogram (as used below in Section IV), which corresponds to  $M \times N \equiv 151 \times 500$  pixels (see Fig. 1). For these map dimensions, and assuming  $t_{\min} = 100$  s, there are  $\approx 2 \times 10^{13}$  possible combinations of ordered pairs making an exhaustive search unfeasible (see Appendix A). However, below we demonstrate that  $T = 2 \times 10^7$  random trials provides sufficient sampling to yield remarkable sensitivity gains with reasonable computational requirements.

Since the *stochtrack* algorithm does not depend on the nearness of seed pixels, it is well-suited for realistic data with instrumental notches (see Fig. 1). It is unaffected by the gaps in  $\rho(t; f)$ .

By design, the *stochtrack* algorithm assumes a particular signal form. Namely, the track is assumed to be reasonably well described by a quadratic Bézier curve with a duration of at least  $t_{\min}$ . (This family of signals includes as a subset all monochromatic tracks with duration of at least  $t_{\min}$ .) In reality, however, the quadratic Bézier curve will be only an approximate fit for an arbitrary monotonic curve. Broadband signals and non-monotonic signals may be poorly fit.

## IV. COMPARISON

In order to demonstrate the *stochtrack* algorithm and compare it to seed-based clustering algorithms we perform a Monte Carlo study. First, we generate Gaussian detector noise following the design sensitivity of Advanced LIGO (aLIGO) at high-power and zero-detuning [17]. Using this simulated noise, we construct spectrograms of  $\rho(t; f)$ . We analyze each spectrogram

with three clustering algorithms: a linear clustering algorithm called *burstegard* [11], a density-based clustering algorithm called *burstcluster* [9], and *stochtrack*. We run two versions of *stochtrack*: a default version with  $T = 2 \times 10^7$  trials and a computationally more expensive deep-search version with  $2 \times 10^8$  trials denoted “*stochtrack* 10 $\times$ .” By running both the default *stochtrack* and *stochtrack* 10 $\times$ , we investigate how detection distance scales with the number of trials.

For each algorithm, we determine the threshold  $\text{SNR}_{\text{tot}}^{\text{th}}$  corresponding to  $\text{FAP} = 0.1\%$  (see Section II). Once we have obtained the thresholds, we perform additional Monte Carlo studies in which a signal is added to the simulated noise. By looping over a range of source distances, we can vary the signal strength, and determine the  $\text{FAP} = 0.1\%$ ,  $\text{FDP} = 50\%$  detection distance  $d_0$  for each algorithm; see Section II. We consider four toy-model waveforms: two down-chirping accretion-disk instability (ADI) waveforms inspired by [6, 7] and calculated following [18] and two up-chirping fallback accretion (FA) powered waveforms from [1]; see Table I.

In the FA model, a newborn neutron star is spun up through fallback accretion following a supernova [1, 2]. The neutron star undergoes a dynamical or secular instability, which induces a time varying quadrupole moment, which in turn powers the emission of narrowband GWs until a black hole is formed and the signal is cut off. In the ADI model, clumps form in the accretion disk surrounding a black hole formed following stellar collapse [6, 7, 18]. The motion of the clumps leads to the emission of narrowband GWs. The ADI waveforms are normalized to assume a GW energy budget of  $E_{\text{GW}} = 0.1M_{\odot}$  [6].

The durations and frequency range of each waveform are given in Table I. The waveform parameters are listed in the Appendix B. The ADI waveforms are analyzed in a band between 100–250 Hz while the FA waveforms are analyzed in a band between 700–1600 Hz.

For our present purposes, we work under the assumption that the GW source location is known, e.g., from an electromagnetic trigger such as a gamma-ray burst or a supernova. We further assume that the time of GW emission is constrained to a small 250 s “on-source” window. While the 250 s window size is comparable to some previous triggered searches for GW bursts, e.g., [19], there are many signal models that would require a significantly larger on-source region [1, 4, 6–8]. Despite this, we restrict the on-source window to 250 s in order to compare different clustering algorithms with a limited computational cost. It is possible to extend this type of analysis to study a larger on-source region at increased computational cost (or with diminished sensitivity at the same computational cost).

We assume that each source is optimally oriented (face-on), which is a reasonable assumption if we assume that the search is carried out following a highly-beamed electromagnetic trigger. We further assume that the detectors are optimally aligned to achieve the maximal possible signal-to-noise ratio. The detection distance ob-

| waveform | duration (s) | $f_{\text{min}}-f_{\text{max}}$ (Hz) | $\delta t \times \delta f$          | $t_{\text{min}}$ |
|----------|--------------|--------------------------------------|-------------------------------------|------------------|
| ADI 1    | 39           | 130–170                              | $1 \text{ s} \times 1 \text{ Hz}$   | 35 s             |
| ADI 2    | 230          | 110–260                              | $1 \text{ s} \times 1 \text{ Hz}$   | 100 s            |
| FA 1     | 25           | 1170–1530                            | $0.5 \text{ s} \times 2 \text{ Hz}$ | 20 s             |
| FA 2     | 200          | 790–1080                             | $1 \text{ s} \times 1 \text{ Hz}$   | 100 s            |

TABLE I: A summary of the waveforms used in our Monte Carlo study. The second and third columns describe the duration and frequency range of the waveform respectively. The fourth column gives the spectrogram resolution used to analyze each waveform. The fifth column specifies the minimum signal duration assumed in each search. The ADI waveforms are down-chirping accretion-disk instability waveforms [6, 7, 18] while the FA waveforms are up-chirping fallback accretion powered waveforms [1, 2].

tained by averaging over detector orientations is  $\approx 60\%$  the value obtained by assuming optimal-aligned distance.

The results are summarized in Table II. We find that, depending on the waveform, the default *stochtrack* improves on the seed-based clustering algorithms by a factor ranging from 150–180% in distance, or equivalently, 320–560% in volume. For *stochtrack* 10 $\times$ , the improvement is 160–200% in distance, or equivalently, 420–740% in volume.

For bright extra-galactic ADI signals with  $E_{\text{GW}} = 0.1M_{\odot}$  [6, 7, 18], we obtain *stochtrack* 10 $\times$  detection distances of  $d_0 = 370\text{--}590$  Mpc. The rate of gamma-ray bursts within this distance range is  $\sim 0.1\text{--}1 \text{ year}^{-1}$ , which suggests that seedless clustering could facilitate the detection of an ADI-type signal by aLIGO [17] / aVirgo [20].

For FA sources [1, 2], we obtain *stochtrack* 10 $\times$  detection distances of  $d_0 = 35\text{--}40$  Mpc [26]. The rate of supernovae in this volume is sufficiently high that aLIGO and aVirgo can expect a promising electromagnetic trigger rate of  $\gtrsim 1 \text{ year}^{-1}$  [1].

The gain in sensitivity is not without added computational cost. On a currently typical computer, the *burstegard* algorithm is capable of analyzing a single  $151 \times 500$  pixel spectrogram in just 1.3 s while the default *stochtrack* algorithm takes 1100 s (18 min) to analyze the same data. The *stochtrack* computation time scales linearly with the number of trials. By increasing the number of trials by a factor of ten, it is possible to increase the detection distance by  $\approx 10\%$ , but the computation time grows to  $1.1 \times 10^4 \text{ s}$  ( $\approx 3 \text{ hr}$ ).

While  $\approx 3 \text{ hr}$  of computing time is not especially burdensome in and of itself, an actual observational analysis will require many ( $\gtrsim 100$ ) pseudo-experiments with time-shifted data. If we further assume that the algorithm is applied to  $\approx 50$  triggers (for example, from gamma-ray bursts), using an on-source region that is larger than the one used here by a factor of  $\approx 50$ , then the estimated computing time is nine weeks on 500 dedicated nodes.

The number of trials can be tuned to match available computational resources. In the event of a detection can-

| waveform | algorithm             | distance |     | volume |
|----------|-----------------------|----------|-----|--------|
|          |                       | absolute | %   | %      |
| ADI 1    | <i>burstcluster</i>   | 330 Mpc  | 90  | 74     |
|          | <i>burstegard</i>     | 370 Mpc  | 100 | 100    |
|          | <i>stochtrack</i>     | 540 Mpc  | 150 | 320    |
|          | <i>stochtrack</i> 10× | 590 Mpc  | 160 | 420    |
| ADI 2    | <i>burstcluster</i>   | 170 Mpc  | 91  | 76     |
|          | <i>burstegard</i>     | 190 Mpc  | 100 | 100    |
|          | <i>stochtrack</i>     | 340 Mpc  | 180 | 560    |
|          | <i>stochtrack</i> 10× | 370 Mpc  | 200 | 740    |
| FA 1     | <i>burstegard</i>     | 17 Mpc   | 100 | 100    |
|          | <i>stochtrack</i>     | 29 Mpc   | 150 | 320    |
|          | <i>stochtrack</i> 10× | 35 Mpc   | 180 | 560    |
| FA 2     | <i>burstegard</i>     | 25 Mpc   | 100 | 100    |
|          | <i>stochtrack</i>     | 36 Mpc   | 150 | 320    |
|          | <i>stochtrack</i> 10× | 40 Mpc   | 160 | 420    |

TABLE II: A comparison of the sensitivity achieved with three different clustering algorithms using aLIGO Monte Carlo noise. *Burstcluster* [9] and *burstegard* [11] use seeds whereas *stochtrack* is seedless. By default, *stochtrack* performs  $T = 2 \times 10^7$  trials. We also report results for *stochtrack* 10× using  $T = 2 \times 10^8$  trials. (Note that *burstcluster* distances are only available for the ADI waveforms since the algorithm is too slow without modification to analyze the larger FA spectrograms.) “Distance” refers to the distance at which a GW source can be detected with false alarm probability = 0.1% and false dismissal probability = 50%. We list both the absolute distance in Mpc and the % relative to the *burstegard* algorithm. The ADI waveforms have been scaled assuming an energy budget of  $E_{\text{GW}} = 0.1 M_{\odot}$ . Volume is given in % relative to the *burstegard* algorithm.

didate, additional trials can be carried out to perform a more sensitive follow-up search. Similarly, a seedless clustering algorithm such as *stochtrack* could be used to follow up on candidates identified by a less sensitive, but computationally cheaper algorithm designed to look for untriggered GW transients in an all-sky, all-time search.

As an additional check, we repeat the comparison of clustering algorithms using initial LIGO noise [27] recolored to match the aLIGO noise curve expected for zero-detuning and high laser power [17]. This allows us to test the performance of the algorithm with non-stationary noise transients and other instrumental artifacts [14, 21]. An unphysical time shift is introduced between the two strain channels in order to remove any coherent signals. The recolored noise results are summarized in Table III. The default *stochtrack* improves on the seed-based clustering algorithms by a factor ranging from 150–180% in distance, or equivalently, 320–560% in volume. For *stochtrack* 10×, the improvement is 160–200% in distance, or equivalently, 420–740% in volume. The similarity between the Monte Carlo and recolored noise results is consistent with previous results [14] and suggests that the expected sensitivity gains from seedless clustering are not dependent on the assumption of

| waveform | algorithm             | distance |     | volume |
|----------|-----------------------|----------|-----|--------|
|          |                       | absolute | %   | %      |
| ADI 1    | <i>burstcluster</i>   | 280 Mpc  | 83  | 57     |
|          | <i>burstegard</i>     | 330 Mpc  | 100 | 100    |
|          | <i>stochtrack</i>     | 540 Mpc  | 160 | 420    |
|          | <i>stochtrack</i> 10× | 540 Mpc  | 160 | 420    |
| ADI 2    | <i>burstcluster</i>   | 159 Mpc  | 91  | 76     |
|          | <i>burstegard</i>     | 170 Mpc  | 100 | 100    |
|          | <i>stochtrack</i>     | 310 Mpc  | 180 | 560    |
|          | <i>stochtrack</i> 10× | 340 Mpc  | 200 | 740    |
| FA 1     | <i>burstegard</i>     | 22 Mpc   | 100 | 100    |
|          | <i>stochtrack</i>     | 32 Mpc   | 150 | 320    |
|          | <i>stochtrack</i> 10× | 35 Mpc   | 160 | 420    |
| FA 2     | <i>burstegard</i>     | 25 Mpc   | 100 | 100    |
|          | <i>stochtrack</i>     | 40 Mpc   | 160 | 420    |
|          | <i>stochtrack</i> 10× | 44 Mpc   | 180 | 560    |

TABLE III: The same as Table II except we use recolored initial LIGO noise with an unphysical time shift instead of Monte Carlo.

idealized detector noise.

## V. CONCLUSIONS

Given a fixed energy budget, a long-lived GW transient produces less excess strain power at any given moment than a short burst. Thus, a long-lived transient is less likely than a short burst with the same total available energy to produce the seed pixels necessary for many traditional clustering algorithms to recover a statistically significant signal. In order to address this, we propose a seedless clustering algorithm called *stochtrack* designed to detect signals too weak to produce seeds. We apply *stochtrack* to several long-lived narrowband signal models and find that it significantly improves detectability compared to two benchmark clustering algorithms, both of which use seeds.

There are a number of ways in which it might be possible to improve the *stochtrack* algorithm. In our current implementation, tracks are fit approximately with quadratic Bézier curves. It may be possible to achieve further improvements in sensitivity using a different, more flexible curve parameterization. The trick with any new parameterization is to better fit test waveforms without expanding the parameter space to the point where the increase in background offsets the gain in signal.

The algorithm may also benefit from improvements in computational efficiency. A more efficient design and/or implementation might reduce the time required to analyze a spectrogram. Reduced computation time, in turn, could facilitate deeper searches (with more trials) and/or searches with large on-source regions. For example, it might be possible to replace the random track generation step with a deterministic process, which more intel-

lightly samples the space of possible curves. One can even imagine the creation of a template bank of curves analogous to the matched filter template banks used for compact binary coalescence searches. (Unlike a matched filter template bank, a *stochtrack* template bank would not contain phase information.)

An area of future research is the application of seedless clustering algorithms to the recovery of compact binary coalescence signals. Of particular interest are regions of parameter space for which it is difficult to create matched filter template banks, e.g., systems with spin and/or eccentricity.

Cornish and Romano have recently emphasized the connection between data analysis algorithms and the signal model for which they are optimal [22]. Following the logic of [23] and [22], *stochtrack* is an optimal search algorithm (in the limit that  $T \rightarrow \infty$ ) for the class of signals described by quadratic Bézier curves in spectrograms of GW power with durations greater than  $t_{\min}$ . Given additional information about the signal model, a seedless clustering algorithm such as *stochtrack* could be tuned appropriately to be more nearly optimal.

### Acknowledgments

We thank Anthony Piro for sharing the fallback accretion waveforms used in this analysis. We thank Shivaraj Kandhasamy and Nelson Christensen for helpful comments. ET is a member of the LIGO Laboratory, supported by funding from United States National Science Foundation. LIGO was constructed by the California Institute of Technology and Massachusetts Institute of Technology with funding from the National Science Foundation and operates under cooperative agreement PHY-0757058. MC is supported by the Winston Churchill Foundation of the United States. This paper has been assigned LIGO document number ligo-p1300103.

### Appendix A: Computational scaling

We estimate the number of possible quadratic Bézier tracks with duration greater than  $t_{\min}$  in a  $M \times N$  spectrogram. For the sake of simplicity, we assume that  $t_{\min}$  is in units of time bins. The number of frequency triplets is given by

$$2 \int_0^M df_3 \int_0^{f_3} df_1 \int_{f_1}^{f_3} df_2 = \frac{M^3}{3}. \quad (\text{A1})$$

Here  $f_1$  is the start frequency,  $f_2$  is the mid frequency, and  $f_3$  is the end frequency. The factor of 2 comes from the fact that the signal can be both up-chirping or down-chirping.

The number of time triplets is given by

$$\int_{t_{\min}}^N dt_3 \int_0^{t_3-t_{\min}} dt_1 \int_{t_1}^{t_3} dt_2 = \frac{N^3}{6} - \frac{t_{\min}^2 N}{2} + \frac{t_{\min}^3}{3}. \quad (\text{A2})$$

Here  $t_1$  is the start time,  $t_2$  is the mid time, and  $t_3$  is the end time. Thus, the total number of possible tracks is

$$\frac{M^3}{3} \left( \frac{N^3}{6} - \frac{t_{\min}^2 N}{2} + \frac{t_{\min}^3}{3} \right). \quad (\text{A3})$$

### Appendix B: Model parameters

The FA waveforms [1, 2] are parameterized by the initial protoneutron star mass  $M_0$ , the maximum neutron star mass  $M_{\max}$ , a dimensionless factor characterizing the supernovae explosion energy  $\eta \approx 0.1$ –10, and the protoneutron star radius  $R_0$ . The two FA waveforms used here assume the following parameters:

| waveform | $M_0$ ( $M_\odot$ ) | $M_{\max}$ ( $M_\odot$ ) | $\eta$ | $R_0$ (km) |
|----------|---------------------|--------------------------|--------|------------|
| FA 1     | 1.3                 | 2.5                      | 10     | 20         |
| FA 2     | 1.3                 | 2.5                      | 1      | 25         |

TABLE IV: Parameters for FA waveforms. See [1] for additional details.

The ADI waveforms [18] are parameterized by black hole mass  $M_{\text{BH}}$ , dimensionless spin parameter  $\alpha^* = [0, 1)$ , the fraction of the accretion disk mass that forms clumps  $\epsilon \approx 0.01$ –0.2, and the torus mass  $m$ . The two ADI waveforms used here assume the following parameters:

| waveform | $M_{\text{BH}}$ ( $M_\odot$ ) | $\alpha$ | $\epsilon$ | $m$ ( $M_\odot$ ) |
|----------|-------------------------------|----------|------------|-------------------|
| ADI 1    | 5                             | 0.3      | 0.05       | 1.5               |
| ADI 2    | 10                            | 0.95     | 0.04       | 1.5               |

TABLE V: Parameters for ADI waveforms. See [18] for additional details.

[1] A. L. Piro and E. Thrane, *Astrophys. J.* **761**, 63 (2012).  
[2] A. L. Piro and C. D. Ott, *Astrophys. J.* **736**, 108 (2011).  
[3] A. L. Piro and E. Pfahl, *Astrophys. J.* **658**, 1173 (2007).  
[4] A. Corsi and P. Mészáros, *Astrophys. J.* **702**, 1171 (2009).

[5] K. Kiuchi, M. Shibata, P. J. Montero, and J. A. Font, *Phys. Rev. Lett.* **106**, 251102 (2011).  
[6] M. H. P. M. van Putten, *Phys. Rev. Lett.* **87**, 091101 (2001).  
[7] M. H. P. M. van Putten, *Astrophys. J. Lett.* **684**, 91



- (2008).
- [8] E. Thrane, S. Kandhasamy, C. D. Ott, et al., Phys. Rev. D **83**, 083004 (2011).
  - [9] R. Khan and S. Chatterji, Classical Quantum Gravity **26**, 155009 (2009).
  - [10] P. Raffai et al., Classical Quantum Gravity **24**, S457 (2007).
  - [11] T. Prestegard and E. Thrane, LIGO DCC p. L1200204 (2012), <https://dcc.ligo.org/cgi-bin/DocDB/ShowDocument?docid=93146>.
  - [12] S. Klimenko and G. Mitselmakher, Classical Quantum Gravity **21**, S1819 (2004).
  - [13] J. Sylvestre, Phys. Rev. D **66**, 102004 (2002).
  - [14] T. Prestegard, E. Thrane, et al., Classical Quantum Gravity **28**, 095018 (2012).
  - [15] B. Abbott et al. (LIGO Scientific Collaboration), Astrophys. J. Lett. **701**, 68 (2009).
  - [16] G. Farin, *Curves and Surfaces for CAGD, Fourth Edition: A Practical Guide* (Academic Press, 1996).
  - [17] Harry, G. M. for the LIGO Scientific Collaboration, Classical Quantum Gravity **27**, 084006 (2010).
  - [18] L. Santamaría and C. D. Ott, LIGO DCC p. T1100093 (2011), <https://dcc.ligo.org/LIGO-T1100093-v2/public>.
  - [19] B. P. Abbott et al., Astrophys. J. **715**, 1438 (2010).
  - [20] Acernese, F. for the Virgo Collaboration, Classical Quantum Gravity **23**, S63 (2006).
  - [21] L. Blackburn et al., Classical Quantum Gravity **25**, 184004 (2008).
  - [22] N. J. Cornish and J. D. Romano, Phys. Rev. D **87**, 122003 (2013).
  - [23] W. G. Anderson, P. R. Brady, J. D. E. Creighton, and É. É. Flanagan, Phys. Rev. D **63**, 042003 (2001).
  - [24] D. Murphy, M. Tse, P. Raffai, et al., Phys. Rev. D **87**, 103008 (2013).
  - [25] One exception is [24], which describes a method for targeting quasimonochromatic signals. By assuming monochromaticity, the space of possible waveforms is dramatically reduced, eliminating the need for seeds. A second exception is the “box search” described in [8], which avoids the use seeds by assuming the signal is broadband.
  - [26] The baseline detection distances for the seed-based *burstegard* clustering algorithm presented here are slightly larger compared to the distances stated in [1]. There are some differences between the two calculations; e.g., the spectrograms are different sizes and here we use a more optimal spectrogram resolution. Most of the discrepancy, however, is accounted for by the use here of a higher-sensitivity aLIGO noise curve.
  - [27] The data are taken in between GPS times 822917487–847549782.